

1 A Appendix

2 Include extra information in the appendix. This section will often be part of the supplemental material.
3 Please see the call on the NeurIPS website for links to additional guides on dataset publication.

- 4 1. Submission introducing new datasets must include the following in the supplementary
5 materials:
 - 6 (a) Dataset documentation and intended uses. Recommended documentation frameworks
7 include datasheets for datasets, dataset nutrition labels, data statements for NLP, and
8 accountability frameworks.
 - 9 (b) URL to website/platform where the dataset/benchmark can be viewed and downloaded
10 by the reviewers.
 - 11 (c) URL to Croissant metadata record documenting the dataset/benchmark available for
12 viewing and downloading by the reviewers. You can create your Croissant metadata
13 using e.g. the Python library available here: <https://github.com/mlcommons/croissant>
 - 14 (d) Author statement that they bear all responsibility in case of violation of rights, etc., and
15 confirmation of the data license.
 - 16 (e) Hosting, licensing, and maintenance plan. The choice of hosting platform is yours, as
17 long as you ensure access to the data (possibly through a curated interface) and will
18 provide the necessary maintenance.
- 19 2. To ensure accessibility, the supplementary materials for datasets must include the following:
 - 20 (a) Links to access the dataset and its metadata. This can be hidden upon submission if the
21 dataset is not yet publicly available but must be added in the camera-ready version. In
22 select cases, e.g when the data can only be released at a later date, this can be added
23 afterward. Simulation environments should link to (open source) code repositories.
 - 24 (b) The dataset itself should ideally use an open and widely used data format. Provide a
25 detailed explanation on how the dataset can be read. For simulation environments, use
26 existing frameworks or explain how they can be used.
 - 27 (c) Long-term preservation: It must be clear that the dataset will be available for a long time,
28 either by uploading to a data repository or by explaining how the authors themselves
29 will ensure this.
 - 30 (d) Explicit license: Authors must choose a license, ideally a CC license for datasets, or an
31 open source license for code (e.g. RL environments).
 - 32 (e) Add structured metadata to a dataset's meta-data page using Web standards (like
33 schema.org and DCAT): This allows it to be discovered and organized by anyone. If
34 you use an existing data repository, this is often done automatically.
 - 35 (f) Highly recommended: a persistent dereferenceable identifier (e.g. a DOI minted by
36 a data repository or a prefix on identifiers.org) for datasets, or a code repository (e.g.
37 GitHub, GitLab,...) for code. If this is not possible or useful, please explain why.
- 38 3. For benchmarks, the supplementary materials must ensure that all results are easily repro-
39 ducible. Where possible, use a reproducibility framework such as the ML reproducibility
40 checklist, or otherwise guarantee that all results can be easily reproduced, i.e. all necessary
41 datasets, code, and evaluation procedures must be accessible and documented.
- 42 4. For papers introducing best practices in creating or curating datasets and benchmarks, the
43 above supplementary materials are not required.

44 **B Supplementary materials**

45 **B.1 Dataset details**

46 **B.1.1 Dataset documentation and intended uses**

47 Dataset documentation is provided in our public public GitHub repository.¹

48 Additional details on how data can be used and how users can utilize the SHDocs Croissant metadata
49 are provided within the data directory of our public GitHub repository

50 **B.1.2 URL to website/platform for dataset/benchmark**

51 The dataset and benchmark are hosted on OneDrive. Download and usage instructions for our dataset
52 and benchmark are detailed in our public GitHub repository

53 **B.1.3 URL to Croissant metadata**

54 The SHDocs Croissant metadata can be accessed on our public GitHub repository here.

55 Additional details on how data can be used and how users can utilize the SHDocs Croissant metadata
56 are provided within the data directory of our public GitHub repository.

57 **B.1.4 Statement on author responsibility**

58 We, the undersigned authors, hereby declare that we collectively bear full and complete responsibility
59 for the content of the work titled "SHDocs: A dataset, benchmark, and method to efficiently generate
60 high-quality, real-world specular highlight data with near-perfect alignment". This includes—but is
61 not limited to—ensuring that all data, materials, and content included in this work are original or
62 appropriately cited and that we have obtained all necessary permissions and rights for any third-party
63 materials used.

64 We affirm that our work complies with all applicable laws and regulations regarding intellectual
65 property, copyright, and ethical standards. In the event of any dispute or violation of rights, we, the
66 authors, will assume all liability and responsibility, and we will take all necessary steps to resolve any
67 issues that may arise.

68 We acknowledge that The Conference and Workshop on Neural Information Processing Systems
69 (NeurIPS) is not responsible for any content-related issues and that we indemnify and hold NeurIPS
70 harmless against any claims, damages, or losses arising from the publication of this work.

71 Signed,

72 Jovin Wei Jie Leong

73 11 June 2024

74 Ming Di Koa

75 11 June 2024

76 Benjamin Wen Bin Cham

77 11 June 2024

78 Shaun Wei Quan Heng

79 11 June 2024

80 **B.1.5 Hosting, licensing, and maintenance plan**

81 The dataset will be hosted on enterprise OneDrive and Google Drive where the links will both be
82 publicly available.

¹<https://github.com/JovinLeong/SHDocs>

83 The data and code are licensed with The MIT License and the licensing details have been included
84 within the public repository.

85 Dataset maintenance will be carried out by all authors; we will be actively working with the dataset
86 for future projects and will use our findings to perform remediation when necessary. We will monitor
87 and address all dataset issues raised on GitHub to ensure that the dataset remains accessible and
88 usable.

89 Code maintenance will similarly be carried out by the authors based on our independent remediation
90 and issues raised by users on GitHub. However, our code maintenance will only be on a best-effort
91 basis—particularly when resolving issues arising from dependencies, development environments,
92 and infrastructure.

93 **B.2 Accessibility**

94 **B.2.1 Links to access the dataset and metadata**

95 The link to access our data and our metadata are included in our public GitHub repository. Addition-
96 ally, SHDocs can be downloaded from the following links:

97 SHDocs raw data: Microsoft OneDrive or Google Drive

98 SHDocs processed data: Microsoft OneDrive or Google Drive

99 **B.2.2 How the dataset can be read**

100 The dataset can be most conveniently read using Croissant to obtain the dataset records in a standard-
101 ized fashion. Users can download the dataset and use the provided Croissant metadata file to load the
102 dataset records. Detailed instructions are provided within the data directory of our public GitHub
103 repository

104 Alternatively, users can manually unzip the dataset and access the data directly. The dataset consists
105 solely of .PNG images and JSON files; these can easily be read by standard libraries available in most
106 programming languages.

107 **B.2.3 Long-term preservation**

108 The dataset will be hosted on enterprise OneDrive and Google Drive where the links will both be
109 publicly available.

110 Dataset maintenance will be carried out by all authors; we will be actively working with the dataset
111 for future projects and will use our findings to perform remediation when necessary. We will monitor
112 and address all dataset issues raised on GitHub to ensure that the dataset remains accessible and
113 usable.

114 Additionally, we will explore the use of data repositories such as Hugging Face Hub later on to
115 maximize data preservation.

116 **B.2.4 Explicit license**

117 The data and code are licensed with The MIT License and the licensing details have been included
118 within the public repository.

119 **B.2.5 Structured metadata**

120 The SHDocs Croissant metadata can be accessed on our public GitHub repository here. Croissant's
121 metadata structure is based on schema.org and is thus compliant to Web standards.

122 **B.2.6 Dereferenceable identifier**

123 The SHDocs uses a public GitHub repository whose unique URL serves as its persistent identifier.

124 **B.3 Reproducibility**

125 To ensure the reproducibility of our benchmark results, we have included everything needed to
126 replicate our findings. This includes our trained model, which is part of the supplementary materi-
127 als. We've also made all the code used in our experiments publicly accessible through our public
128 GitHub repository. The repository contains all necessary scripts, configuration files, and dependency
129 information.

130 We've provided detailed documentation within the repository to guide you through the entire process.
131 This documentation covers how to set up the environment, run the code, and follow the exact steps
132 we used for model training, evaluation, and result generation.

133 All external models and datasets by other authors that we used in our benchmark are clearly cited and
134 publicly available.

135 Where possible, we have followed The Machine Learning Reproducibility Checklist v2.0 to ensure
136 our experiments are documented and reproducible. By providing all these resources, we aim to make
137 it straightforward for anyone to verify and build upon our work.